



MSiA Seminar Series

# From Anomaly Detection to Data Visualization: In the Trenches of the CTIO's Office

Christopher Laporte

February 26, 2019



**Jet Propulsion Laboratory**  
California Institute of Technology

Copyright 2019 California Institute of Technology. U.S. Government  
sponsorship acknowledged

- **JPL**
- **Data Science @ JPL**
- **Case Projects**
- **Wrap-up**
- **Q/A**



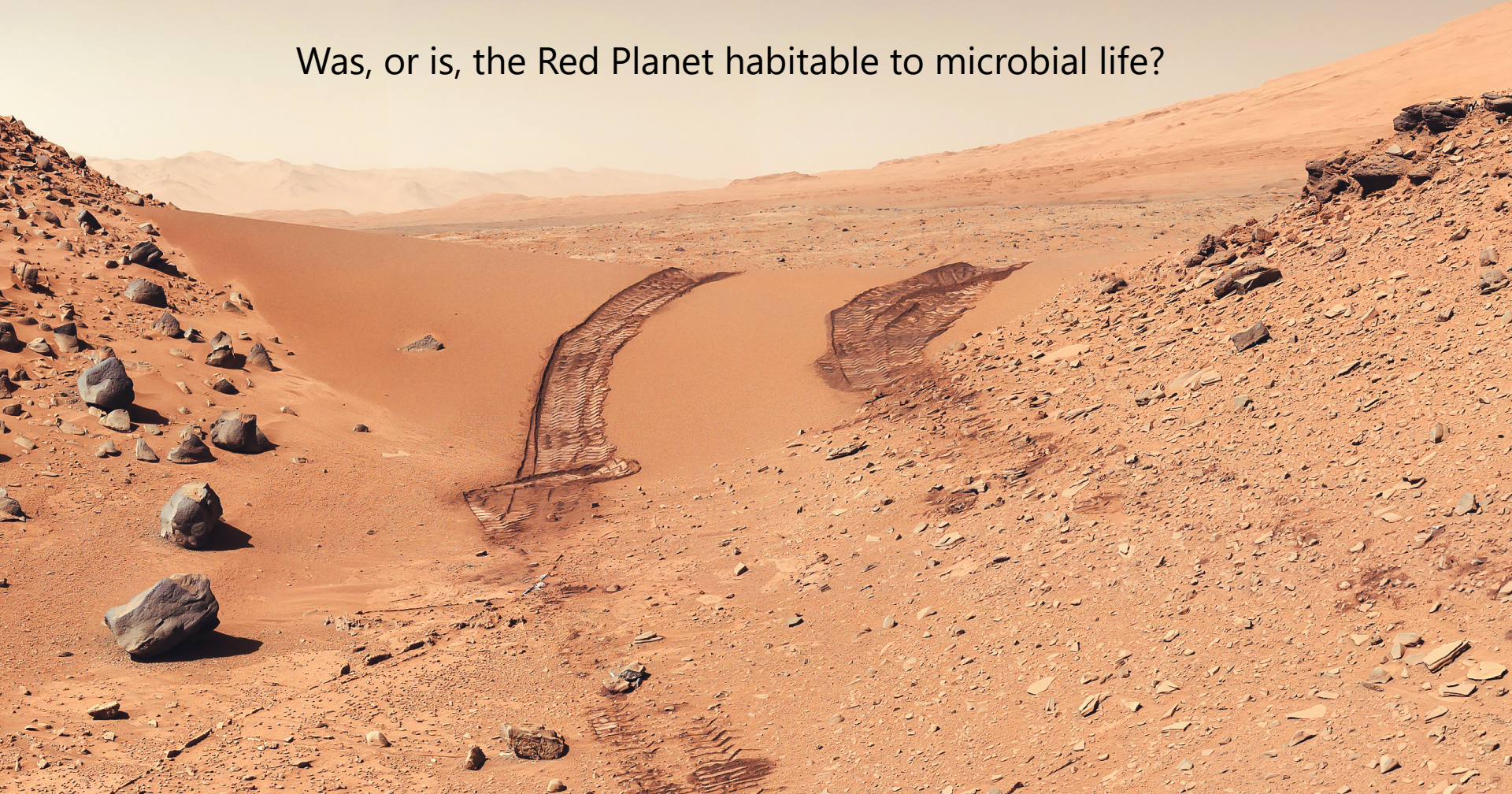
# 1. JPL



What can Jupiter's formation and evolution  
tell us about our solar system?



Was, or is, the Red Planet habitable to microbial life?



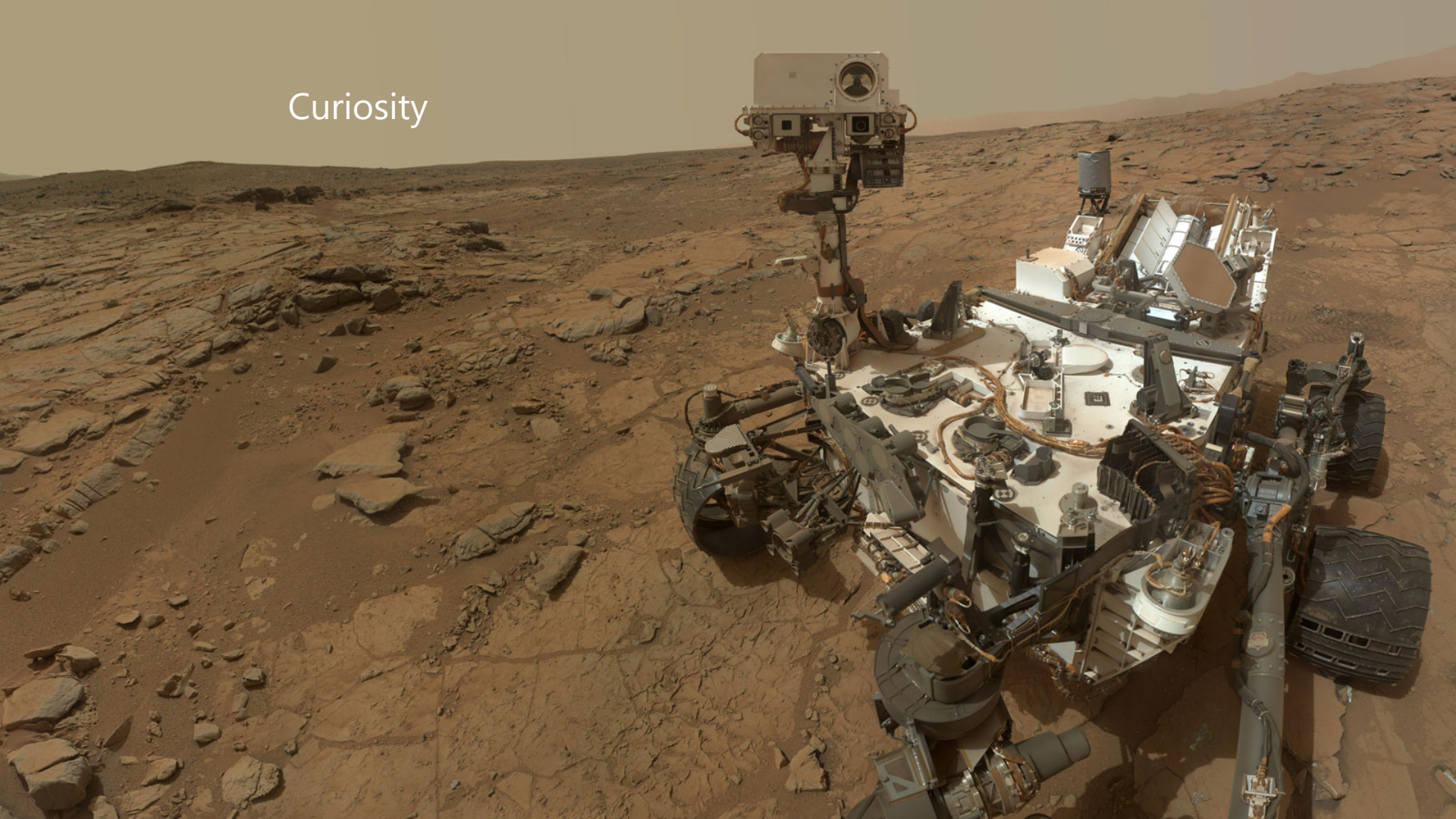




Juno

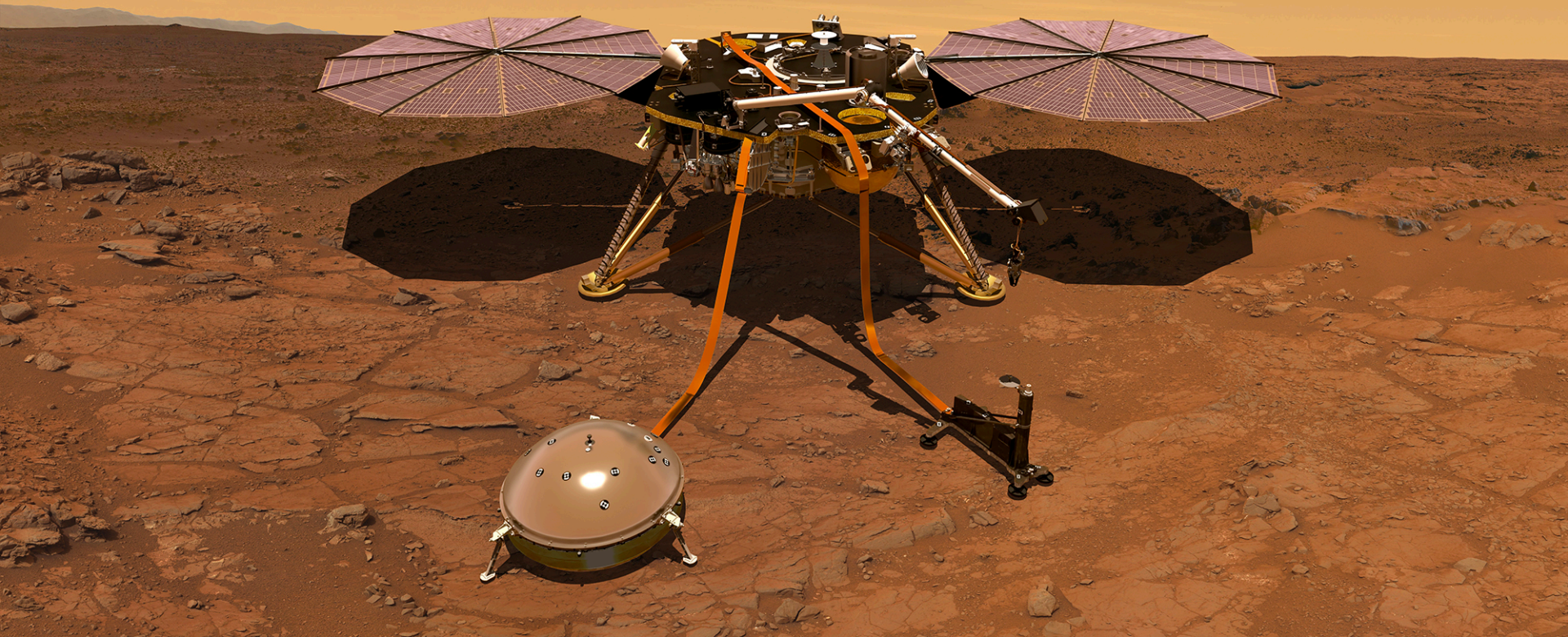


Curiosity





Insight





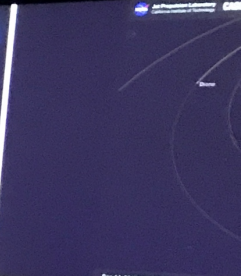
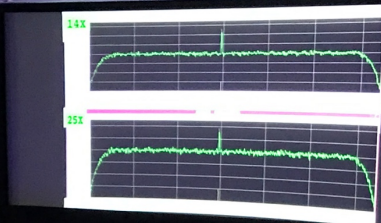
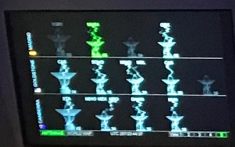








Laboratory  
of Technology



CHARLES ELACHI MISSION CONTROL CENTER

MISSION CONTROL

DATA CONTROLLER

DEEP SPACE NETWORK

DEEP SPACE NETWORK

DEEP SPACE NETWORK

Charles Elachi Mission Control Center



# Snapshot

- Over 45 current (in-flight) missions
- Over 15 future missions currently in development
- Operates the Deep Space Network (DSN)
- Scope and volume of scientific data is large (NISAR will produce 3-5 TB daily)



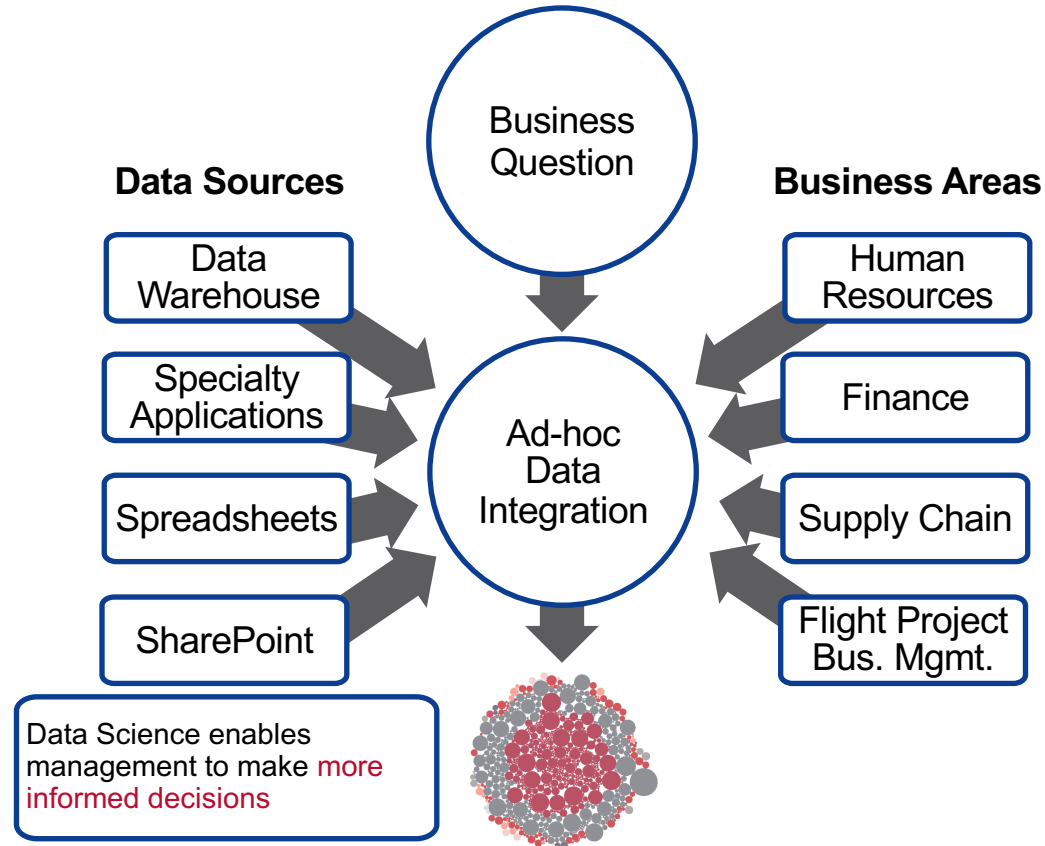
A single DSN 70-meter radio antenna in Goldstone, CA

# How can we as data scientists aid JPL in its mission?

## 2. Data Science @ JPL



# Business IT Data Science

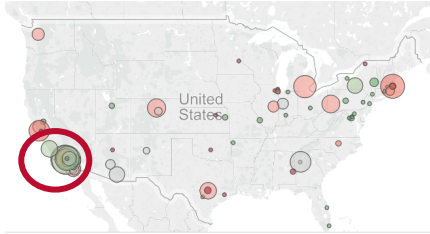


Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California

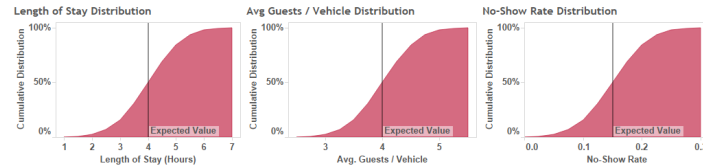
Institute of Technology.

# Business IT Data Science: Sample Projects

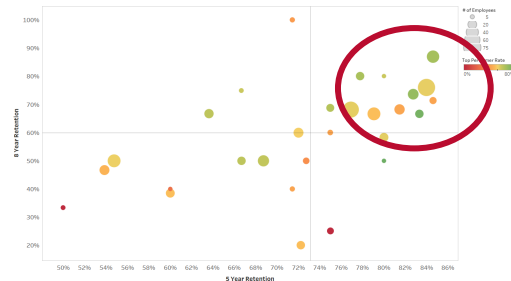
Q: How can we **improve Early Career Hire retention?**



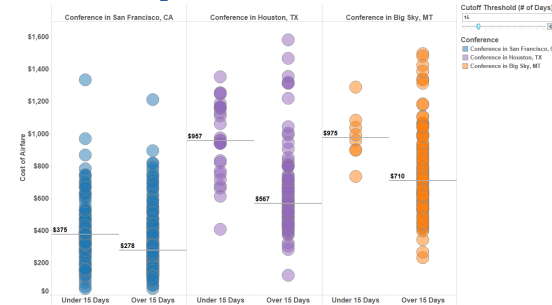
Q: How many tickets should we release for Explore JPL and how should we **allocate the tickets** throughout the day?



Q: Which **schools** provide us with the **most successful employees?**



Q: How can we reduce **conference travel expenses?**



Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.

# Chief Technology and Innovation Office

- Collection of data scientists, cloud engineers, software developers, and data visualization gurus
- Infuse new technologies and techniques into the way we do things at JPL
- Automated intelligence, digital transformation, unstructured information management, open source, cyber security, chatbots, IoT, next-gen robotics and flight hardware



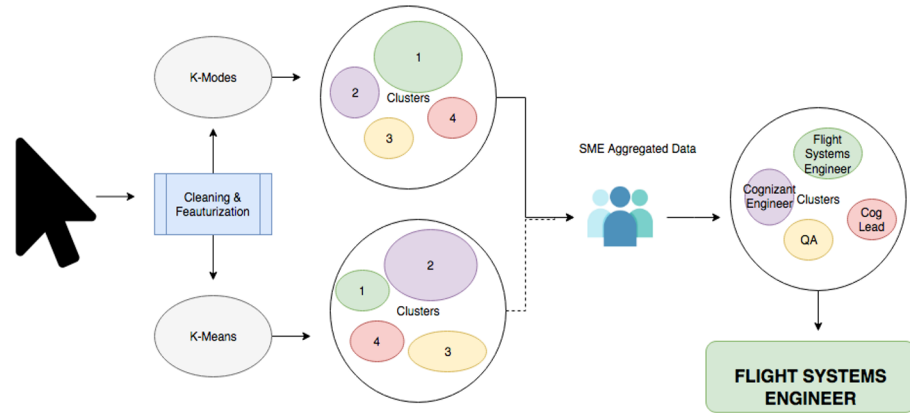
# How We Work

- Talk to everyone and build trust
- Identify and seize moments of engagement with passionate end-users
- Rapidly prototype and iterate
- Focus on the user experience



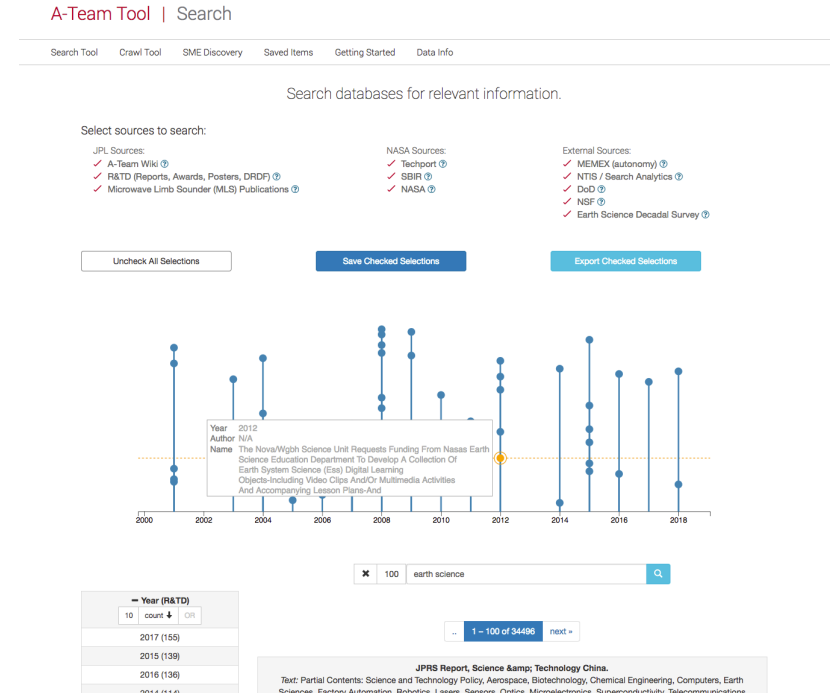
# Engineering Data Management Initiative

- What do our employees do?
- Classify JPL Roles using clickstream data from a variety of tools



# Foundry Data Science

- A-team studies: Early stage mission formulation and feasibility assessment
- Need to know who to include (SMEs), what historical information is available
- A-team tool



# Problem Reporting System

- Henosis: A python framework for deploying recommendation models for form fields
  - Open Issues! <https://github.com/vc1492a/henosis/issues>
- Identifying minority class labels from limited training data
  - Ex: spacecraft safings, escapes

# JPL Open Source Rover

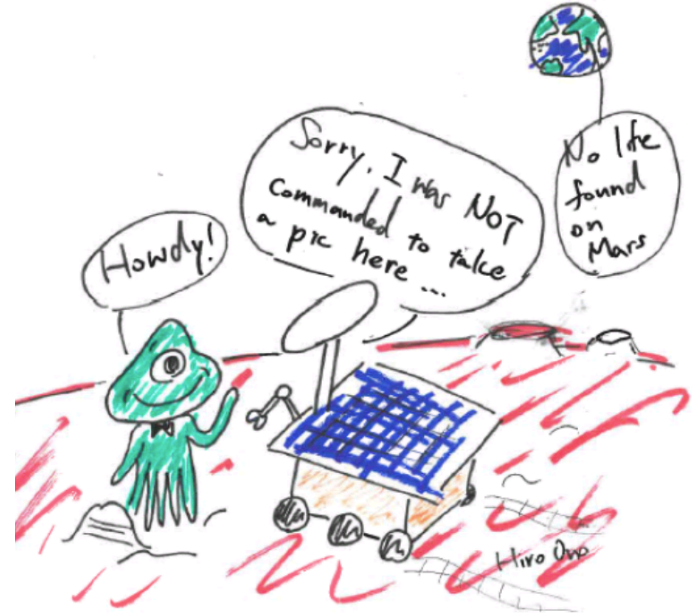
- DIY Rover you can build at home
- <https://opensourcerover.jpl.nasa.gov/>





# Rover Drive-By Science

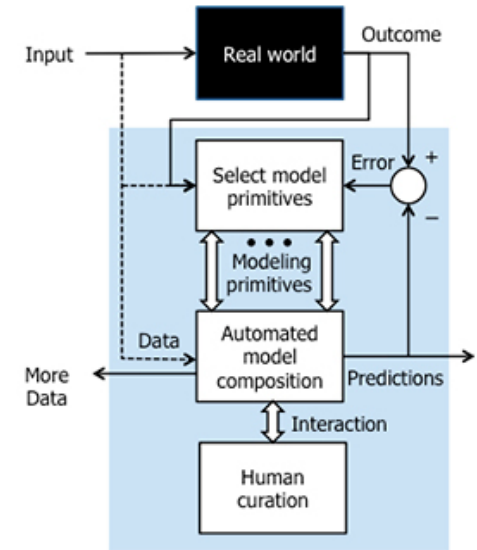
- Current major bottleneck for AI on rovers: extremely limited on-board computation resources
- What would high performance spacecraft computing enable for future missions?



Comic by PI Hiro Ono detailing the "Unnoticed Green Monster Problem" (UGMP)

# Data-Driven Discovery of Models (DARPA)

- Automate the development of machine learning pipelines
- Allow SMEs to analyze data without the need for a data scientist
- Architecting and implementing a library of ML primitives
- Facilitating and cooperation and collaboration between the 23 performers

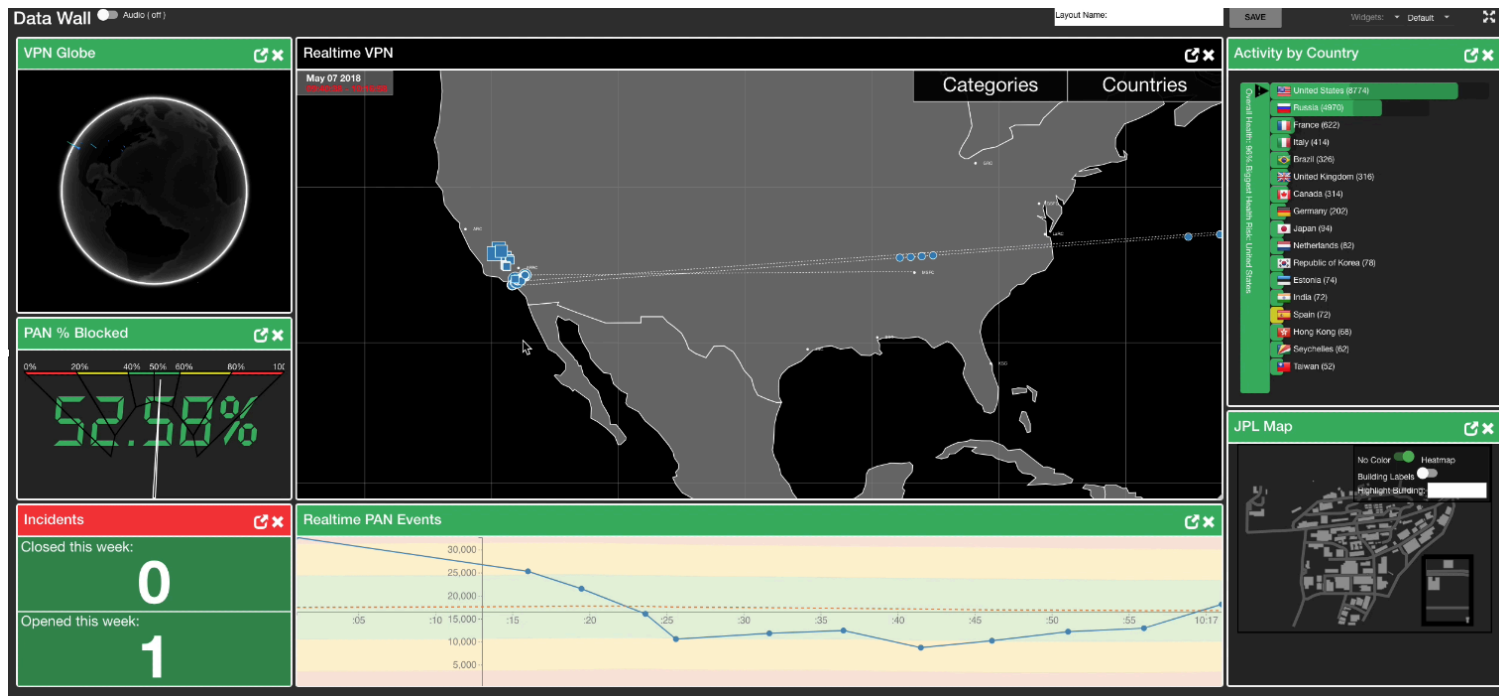




# Active Social Engineering Defense (DARPA)

- Proactively respond to social engineering attacks (eg, phishing emails) with chatbots
- Identify and engage attacks, eventually turning them over to law enforcement
- Develop a test environment that will utilize JPL's email system and evaluation methodologies for the performers

# Cyber Security Data Wall





## And much more...

- Small satellite data science
- Cloud engineering
- Next-gen flight hardware/robotics
- MGSS (multi-mission ground systems and services) open source policy
- Intelligent assistants/chatbots

### **3. Case Project: Cyber Security Visualization**

# Snapshot

- ~ 330,000 events are observed by our firewall... every minute... after filtering
- Use a host of 3<sup>rd</sup> party software/tools to help monitor network
- Even after best filtering attempts, security engineers are still left with ~650 potential threats to review / day
- Cost of a single successful attack could be astronomical



New Search

Save As

Close

All time

```

sourcetype=fgt_traffic src=192.168.225.* NOT (dest=192.168.* OR dest=10.* OR dest=8.8.4.4 OR dest=8.8.8.8 OR dest=224.*)
| eventsstats sum(bytes_out) AS total_bytes_out by src
| stats sum(bytes_in) as bytes_in sum(bytes_out) values(total_bytes_out) AS total_bytes_out by src dest
| eval percent_bytes_out = bytes_out/total_bytes_out * 100
| table src dest bytes_in bytes_out total_bytes_out percent_bytes_out
| where percent_bytes_out > 60
| sort - percent_bytes_out dest

```

0 events (before 6/20/17 1:43:13.000 AM) No Event Sampling

Job

Smart Mode

Events

Patterns

Statistics (31)

Visualization

20 Per Page

Format

Preview

< Prev

1

2

Next >

src	dest	bytes_in	bytes_out	total_bytes_out	percent_bytes_out
192.168.225.80	185.151.160.15	356396	395942	512060	77.323361
192.168.225.110	185.151.160.15	178198	197971	281528	70.320181
192.168.225.166	185.151.160.15	178198	197971	284473	69.592193
192.168.225.121	185.151.160.15	178198	197971	286009	69.218451
192.168.225.142	185.151.160.15	178198	197971	292279	67.733570
192.168.225.66	185.151.160.15	178198	197971	297210	66.609805
192.168.225.79	185.151.160.15	178198	197971	299518	66.096528
192.168.225.39	185.151.160.15	178198	197971	300800	65.814827
192.168.225.200	185.151.160.15	178198	197971	301793	65.598274
192.168.225.23	185.151.160.15	178198	197971	303019	65.332867
192.168.225.54	185.151.160.15	178198	197971	303865	65.150972
192.168.225.143	185.151.160.15	178198	197971	303975	65.127395
192.168.225.88	185.151.160.15	178198	197971	305986	64.699365
192.168.225.201	185.151.160.15	178198	197971	306810	64.525602

Splunk DB Connect

Search

Database Info

Database Query

Searches

Settings

Splunk DB Connect v1

Edit: sample\_h

+ Add Panel

+ Add Input

+ Edit Source

Done

Untitled

Q

BUSINESS_DATE	COMPANY_CODE	LINE_NO	ORACLE_GL	CONTRACT_CODE	APL_CODE	CUSTOMER_ID	CURRENCY	BALANCE	CCY_BALANCE	PROD_CODE	BRANCH	ASST_TYPE
1391106600.000	EH0010001	EHGTB.5342	311133	1133066615	AC	3073609	ETB	280	280	1020	4001	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133065856	AC	3073508	ETB	2080.5	2080.5	1020	4026	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133064612	AC	3073377	ETB	6499	6499	1020	4001	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133063438	AC	3073267	ETB	25	25	1020	4095	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133063373	AC	3073260	ETB	420	420	1020	4055	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133063365	AC	3072157	ETB	62	62	1020	4008	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133063292	AC	3073231	ETB	3960.91	3960.91	1020	4012	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133063071	AC	3073225	ETB	420	420	1020	4055	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133062911	AC	3073177	ETB	66055	66055	1020	4509	CREDIT
1391106600.000	EH0010001	EHGTB.5342	311133	1133244432	AC	3091000	ETB	628.28	628.28	1020	4154	CREDIT

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.

# DEMO

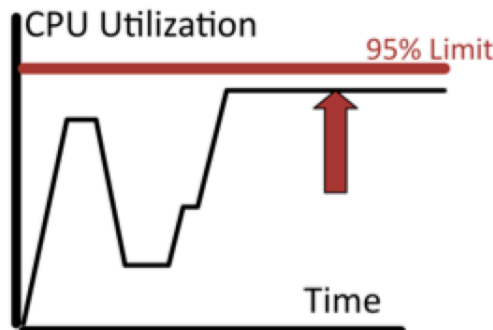
### **3. Case Project: Telemetry Anomaly Detection**



# Motivation

- Thresholding, expert systems

- Reliance on expert knowledge
- Custom
- Not complete
- Accuracy
- Appropriate limits change



Simple example of anomaly that would be undetected by a threshold

~40% of anomalies in experiments are of this nature

- Increasing data rates

- SWOT, NISAR = 3-5 TB daily

- Smaller missions (e.g. cubesats)

- Less people for ops

- High volumes of testbed data

- Investigative aspect

- Focused, prioritized telemetry review
- Help with causal fault analysis
  - What anomalies were detected leading up to a failure?

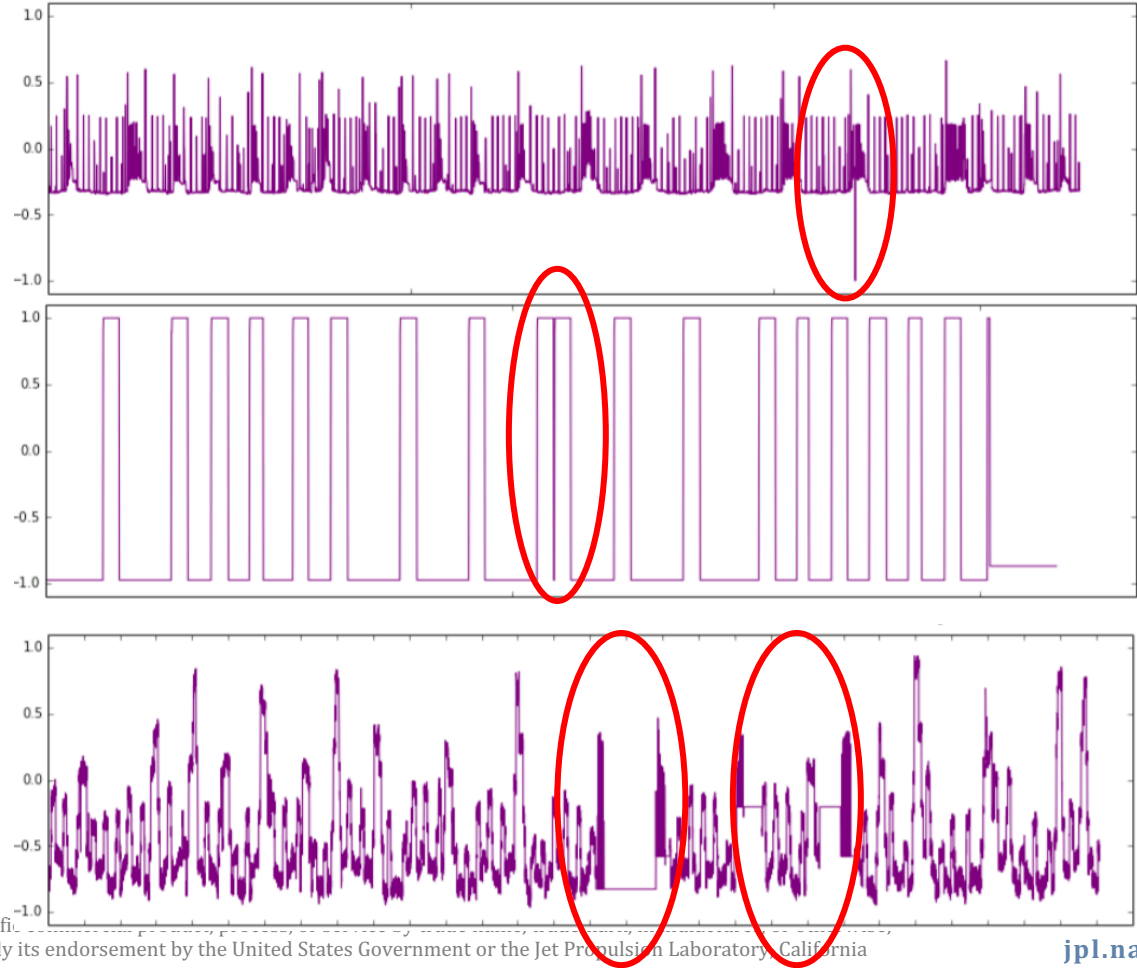
# Anomaly Categories

Chandola et al. 2007

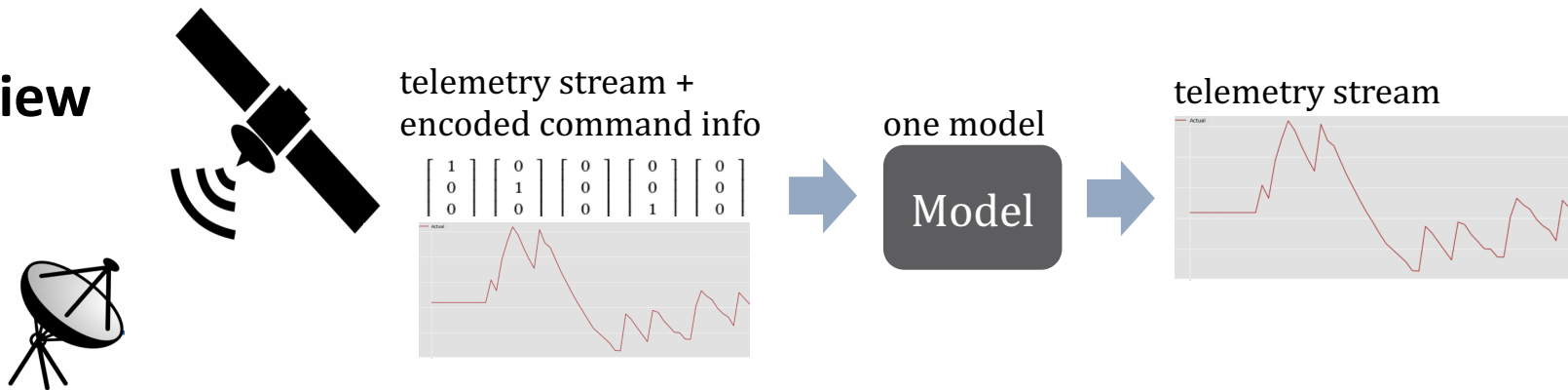
## 1. Point

## 2. Contextual

## 3. Collective (sequential)



# Overview



- Use Recurrent Neural Networks (LSTMs) to predict incoming telemetry values using recent telemetry, commands, and event records (EVRs) as inputs
- Where predictions are substantially different from actual telemetry values, these are identified as potentially anomalous events
  - Novel method for defining “substantially different”
- <https://www.kdd.org/kdd2018/accepted-papers/view/detecting-spacecraft-anomalies-using-lstms-and-nonparametric-dynamic-thresh>

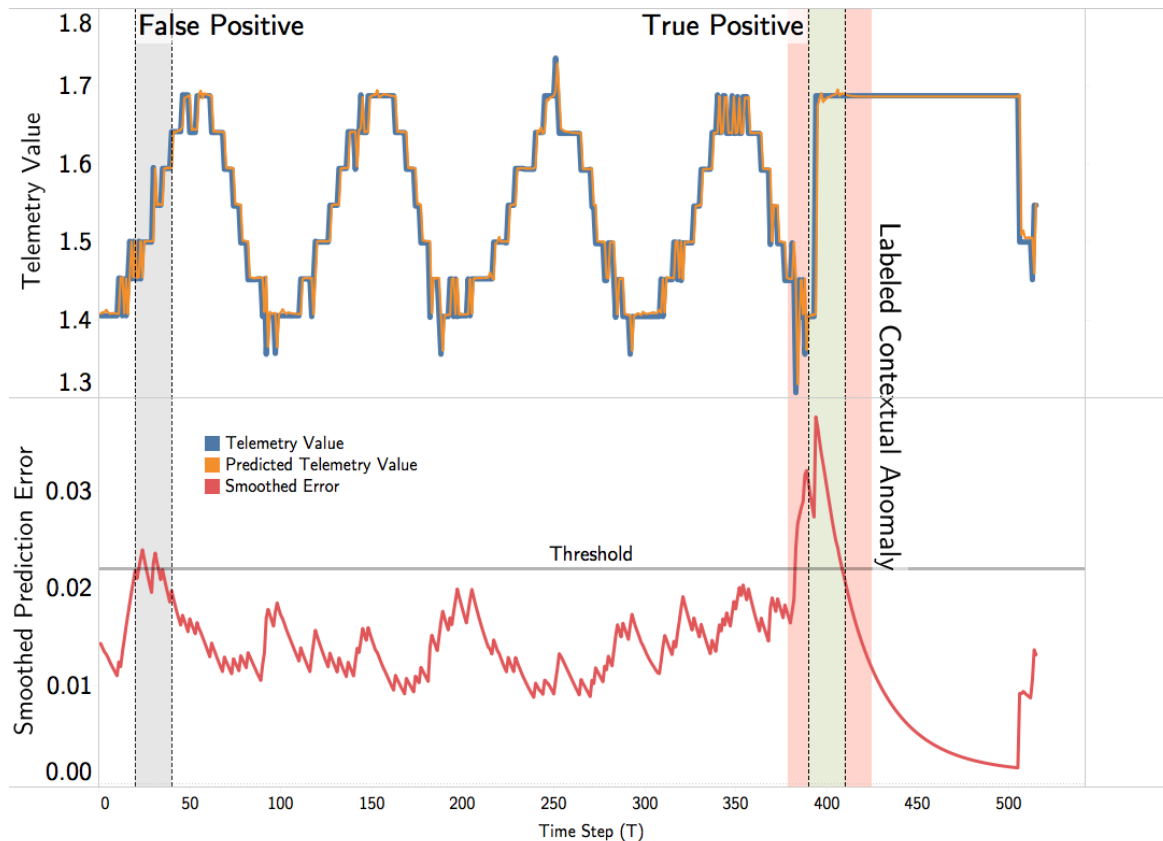


# Single-Channel Prediction



# Reconstruction Errors

Actuals and Prediction

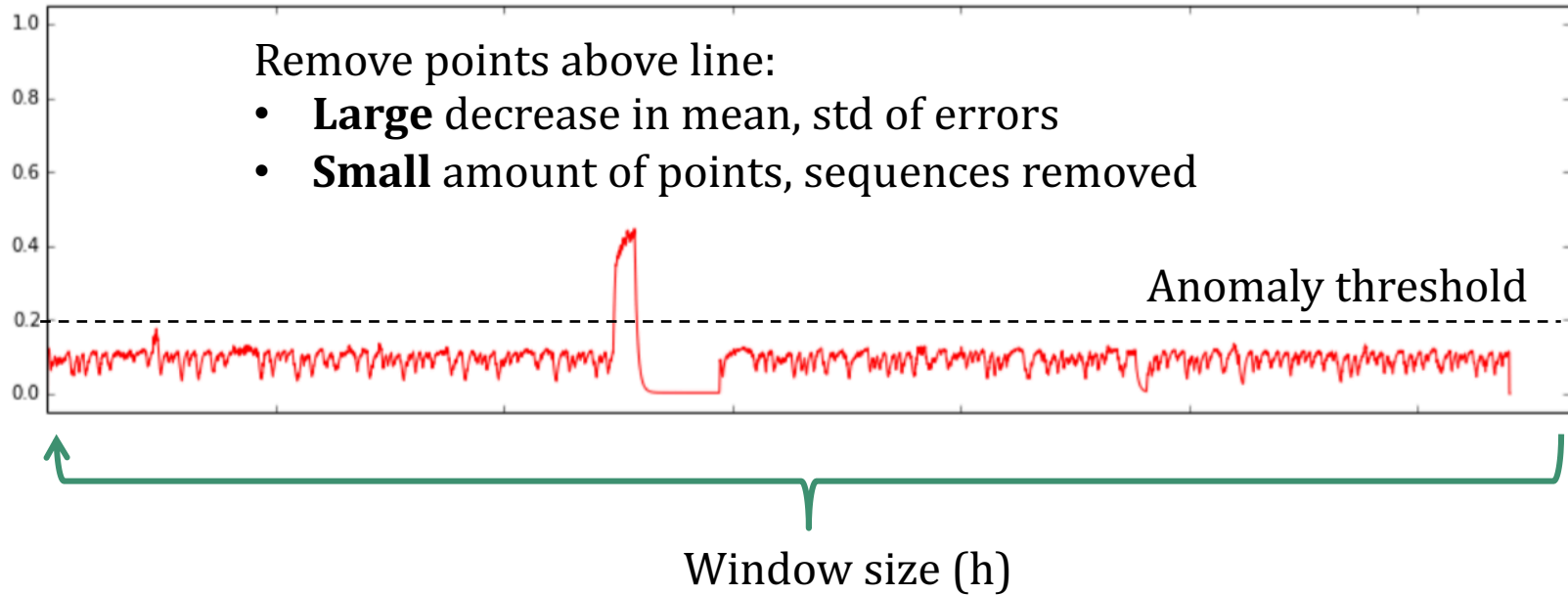


Raw Reconstruction Error

# Dynamic Anomaly Threshold

## Anomalous

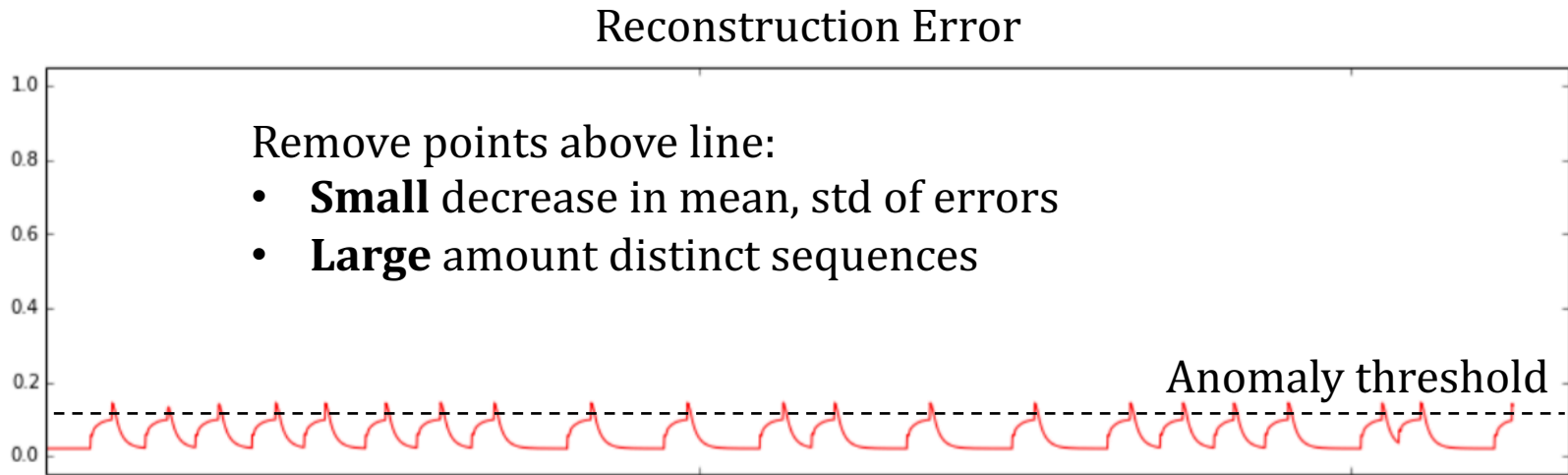
Reconstruction Error





# Dynamic Anomaly Threshold

## Nominal



# Dynamic Anomaly Threshold

Smoothed errors  $\mathbf{e}_s = [e_s^{(t-h)}, \dots, e_s^{(t-l_s)}, \dots, e_s^{(t-1)}, e_s^{(t)}]$

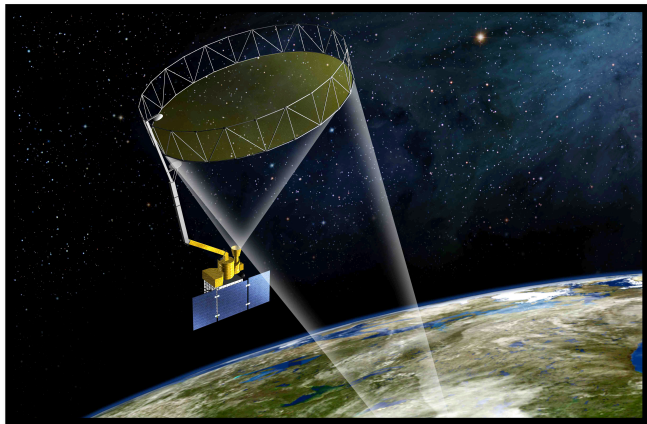
Candidate thresholds  $\epsilon = \mu(\mathbf{e}_s) + \mathbf{z}\sigma(\mathbf{e}_s)$

Threshold  $\epsilon = \operatorname{argmax}(\epsilon) = \frac{\Delta\mu(\mathbf{e}_s)/\mu(\mathbf{e}_s) + (\Delta\sigma(\mathbf{e}_s)/\sigma(\mathbf{e}_s))}{n(\mathbf{e}_a) + n(\mathbf{E}_{seq})^2}$

Definitions

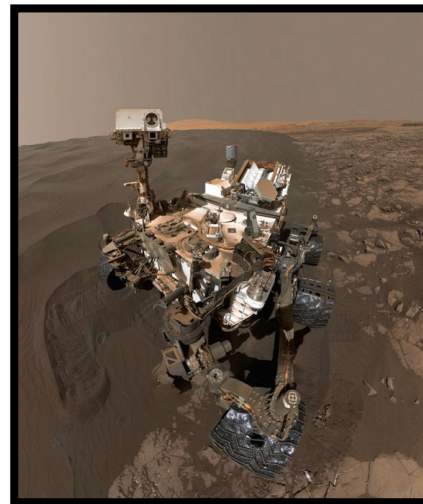
$$\begin{aligned}\Delta\mu(\mathbf{e}_s) &= \mu(\mathbf{e}_s) - \mu(\{e_s \in \mathbf{e}_s | e_s < \epsilon\}) \\ \Delta\sigma(\mathbf{e}_s) &= \sigma(\mathbf{e}_s) - \sigma(\{e_s \in \mathbf{e}_s | e_s < \epsilon\}) \\ \mathbf{e}_a &= \{e_s \in \mathbf{e}_s | e_s > \epsilon\} \\ \mathbf{E}_{seq} &= \text{continuous sequences of } e_a \in \mathbf{e}_a\end{aligned}$$

# Experiments – Two Representative Spacecraft



## Soil Moisture Active Passive (SMAP)

- Higher, more consistent data rates
- Fewer, more routine behaviors

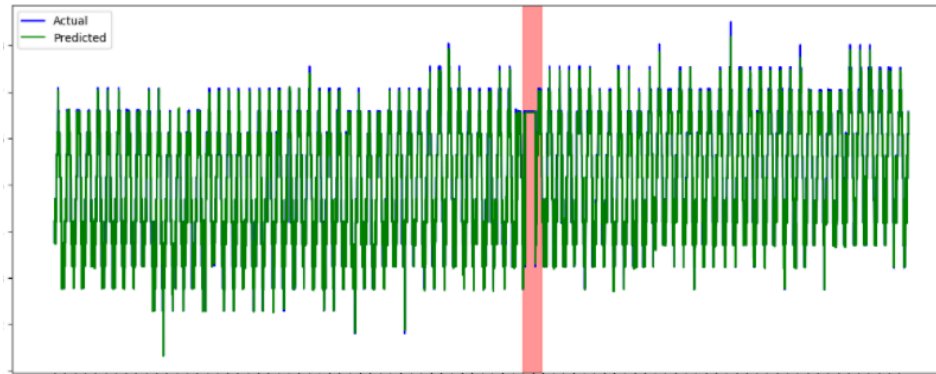


## Mars Science Laboratory (Curiosity or MSL)

- More channels (12k)
- Less data, less consistent delivery
- Extremely varied behaviors
  - Training on recent data isn't enough

# Experiments – Incident Surprise, Anomaly Reports (ISAs)

- Scraped ISAs to find mentions of telemetry channels
  - Ex. “On DOY 192, in the time range from 09:21z through 10:47z, the following channels were found to have odd constant values: A-3, A-4, A-5, A-6, G-3”



- Labeled anomalous ranges for 112 unique ISA anomalies
- Significant portion of contextual anomalies (39%)

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.



# Validation: Predicting ISAs

- Identified all Incident, Surprise, Anomaly (ISA) reports that were apparent in telemetry (EHA) for SMAP and MSL
- Ran Telemanom system over time period surrounding each ISA to see if system would have detected the anomaly



## Results

Thresholding Approach	Precision	Recall	$F_{0.5}$ score
<b>Non-Parametric w/ Pruning (<math>p = 0.13</math>)</b>			
MSL	92.6%	69.4%	<b>0.69</b>
SMAP	85.5%	85.5%	<b>0.71</b>
Total	87.5%	80.0%	<b>0.71</b>

80% of all ISAs  
were identified  
(~115 in total)

	Recall - <i>point</i>	Recall - <i>contextual</i>
MSL	78.9%	58.8%
SMAP	95.3%	76.0%
Total	90.3%	69.0%

*Contextual* anomalies  
are those that are not  
detectable by  
thresholds (0% recall)

# Current Work: MSL

- Extending Telemanom to rovers/planetary missions
  - Prediction of telemetry is harder with more variety and irregularity of behaviors
  - Models need more training and detailed inputs surrounding commands and EVRs
- Early progress
  - Detected Martian sandstorm early with small number of Thermal channels
  - Achieving very high prediction accuracy for thermal channels (~98%)

# Future Work

- Research new methods of dimensionality reduction for our EVR encoding
- Refactor code base to generalized and modular state
- Provide an API and frontend adjustments to allow for the training of multiple channels within a single LSTM
- Research and compare new modeling methods for time sequenced data

## 4. Wrap-up



# Review: How We Work

- Talk to everyone and build trust
- Identify and seize moments of engagement with passionate end-users
- Rapidly prototype and iterate
- Focus on the user experience

# Thanks

## 5. Q/A



**Jet Propulsion Laboratory**  
California Institute of Technology

---

[jpl.nasa.gov](https://jpl.nasa.gov)



# Formulation

Model inputs at step  $t$

$$X = \left\{ \begin{bmatrix} x_1^{(t-l_s)} \\ x_2^{(t-l_s)} \\ \vdots \\ x_m^{(t-l_s)} \end{bmatrix}, \dots, \begin{bmatrix} x_1^{(t-1)} \\ x_2^{(t-1)} \\ \vdots \\ x_m^{(t-1)} \end{bmatrix}, \begin{bmatrix} x_1^{(t)} \\ x_2^{(t)} \\ \vdots \\ x_m^{(t)} \end{bmatrix}, \begin{bmatrix} x_1^{(t+1)} \\ x_2^{(t+1)} \\ \vdots \\ x_m^{(t+1)} = y^{(t)} \end{bmatrix} \right\}$$

Telemetry Values

$h$  = historical window of errors  
 $l_s$  = sequence length

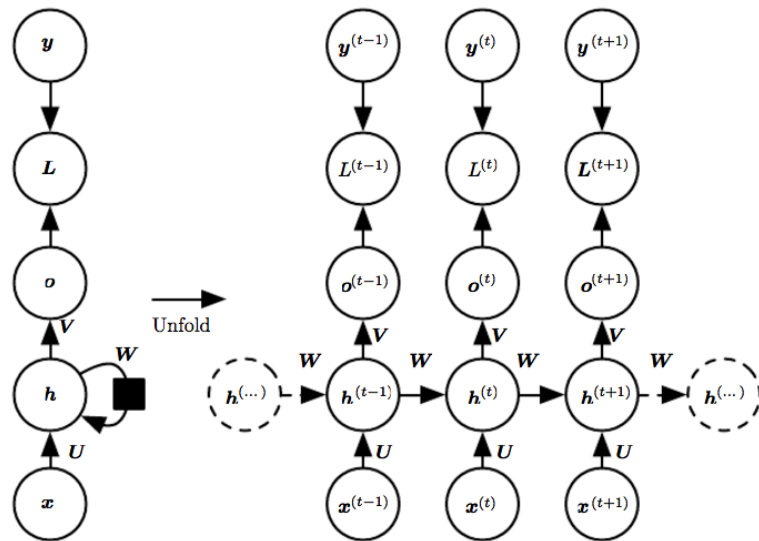
$$e^{(t)} = [\hat{y}^{(t)} - y^{(t)}]$$

$$\mathbf{e} = [e^{(t-h)}, \dots, e^{(t-l_s)}, \dots, e^{(t)}]$$

# Recurrent Neural Nets

- Memory (lossy summary)
- Parameter sharing
  - Extend model to apply to different lengths and generalize across time steps
    - Don't have to have separate parameters for each time value
- Recurrence
  - Always has same input size regardless of sequence length

$$\begin{aligned} \mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}). \end{aligned}$$



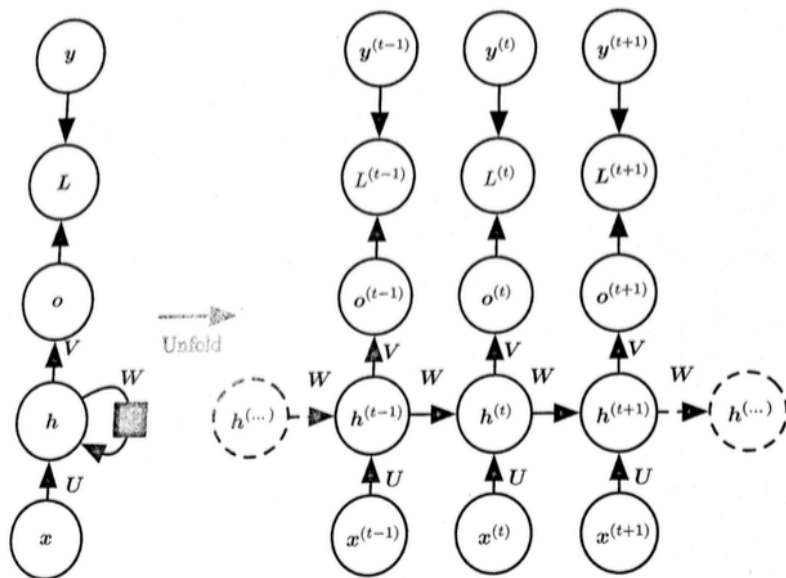
Ian Goodfellow, Yoshua Bengio, Aaron Courville, 2016. *Deep Learning*. MIT Press. <http://deeplearningbook.org>.

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California

Institute of Technology.

# From RNNs to LSTMs (Goodfellow et. al, 2016)

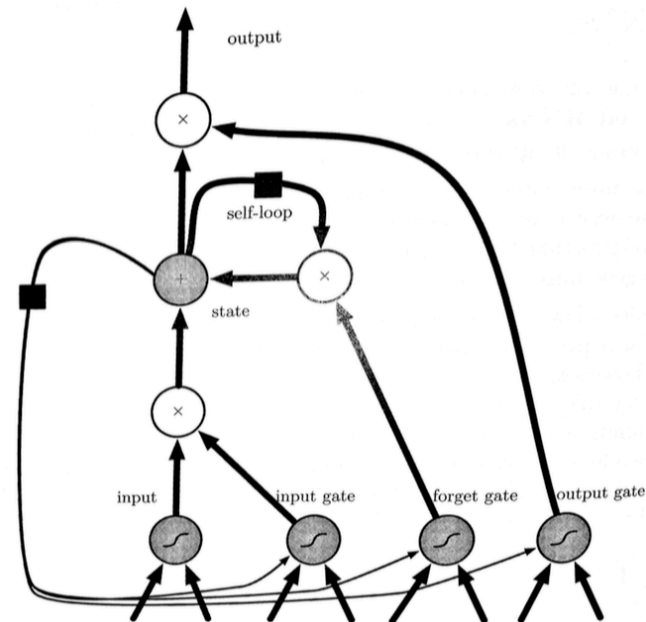
## RNN



Core contribution (1997): Self-loops

Crucial addition (2000): Condition loop on context (with another hidden unit)

## LSTM



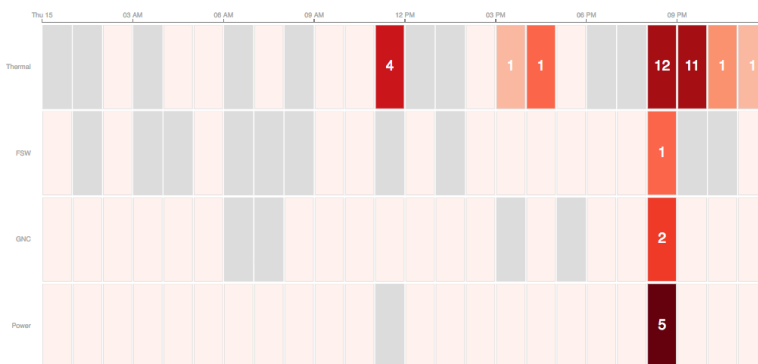
Ian Goodfellow, Yoshua Bengio, Aaron Courville, 2016. *Deep Learning*. MIT Press. <http://deeplearningbook.org>.

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California

Institute of Technology.

# Initial Pilot: SMAP

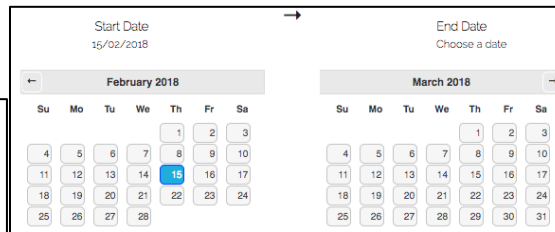
- Deployed end-to-end autonomous system
- Monitored ~750 core telemetry channels from Aug 2017 – May 2018
  - Detected multiple verified anomalous events
    - Partial eclipse (Feb 15, 2018)
- Radar (HPA) failure investigation
  - Ran system ~2 months prior to failure, detected many of same telemetry oddities that were identified during peer review process following failure





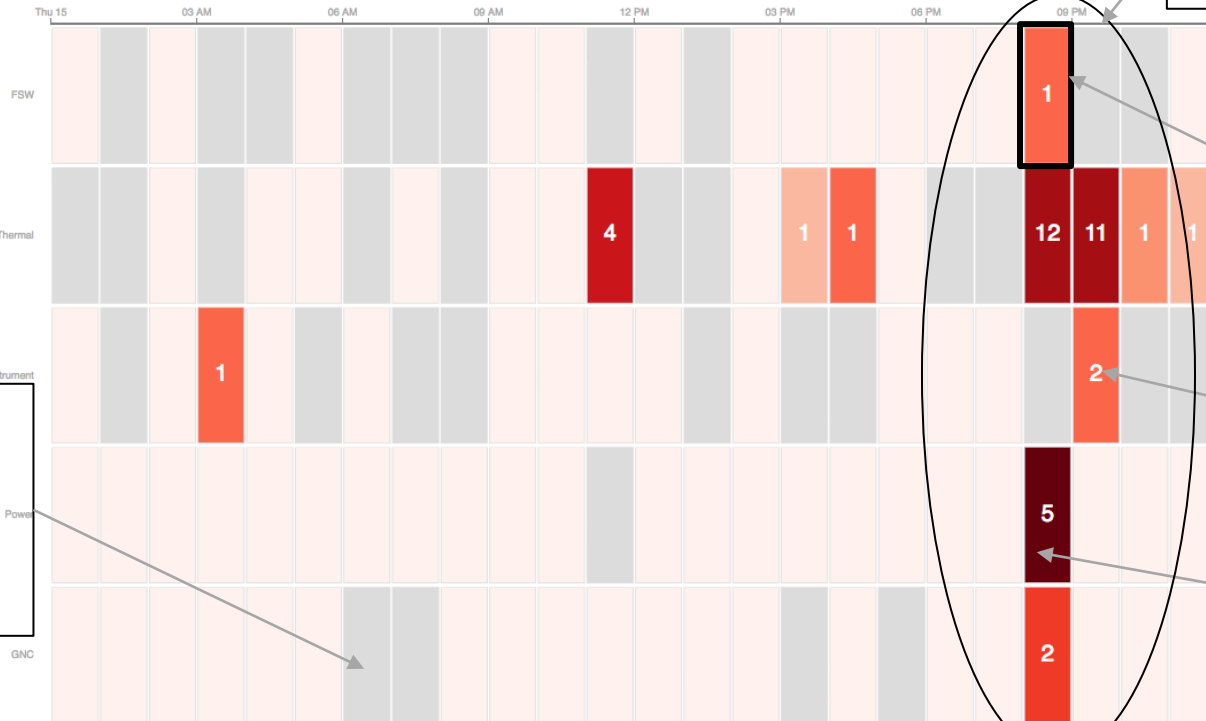
# Interface: Top-Level Summary

Start by selecting a start and end date to look at



February 15<sup>th</sup>, 2018  
Partial Solar Eclipse anomaly

Subsystem



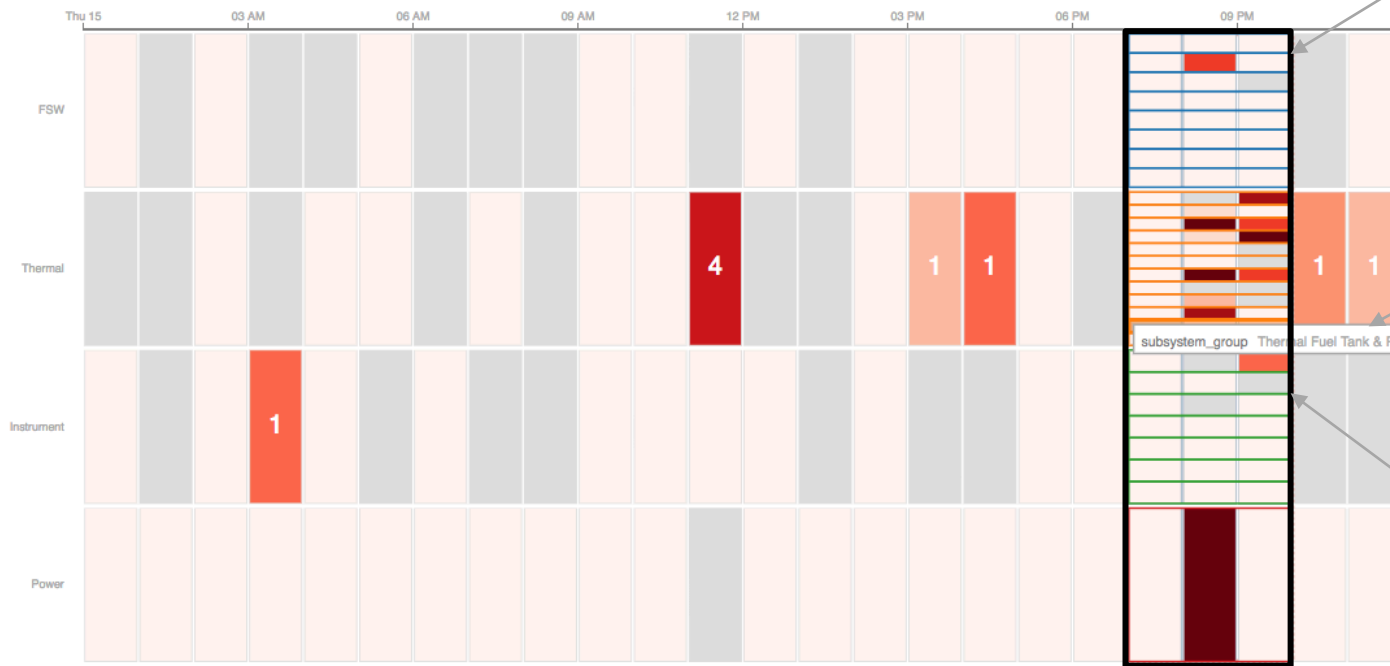
Each box indicates one hour aggregate time window (adjustable)

Count of channel anomalies in subsystem during hour time window

Darker color indicates

Gray boxes are potential anomalies that the system has learned are false positives with high likelihood ("suppressed" anomalies)

# Interface: Drilldown



Clicking and dragging across an area allows for looking down a level to channel groups with subsystems

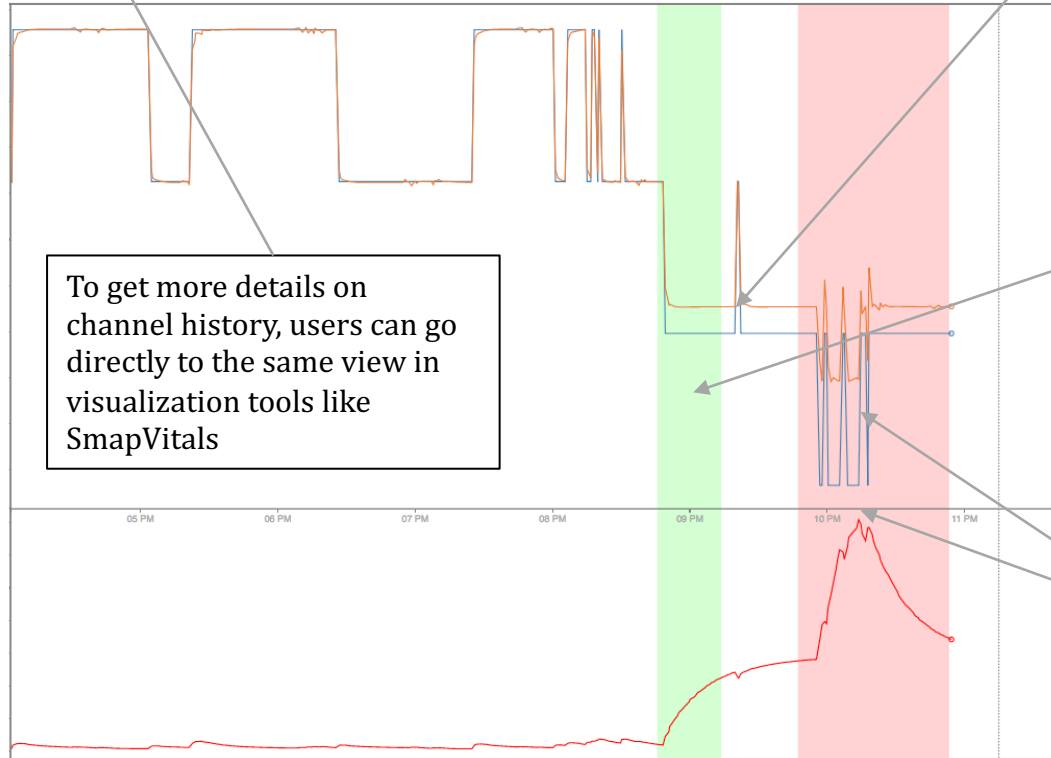
Each row represents a group of channels and hovering shows the group name

Clicking takes the user into a similar view but in the next level down for the selected window

# Interface: Drilldown (cont.)

BACK SMAP VITALS SAVE CHANGES

Anomaly (anom15187563460001518758028000) changed to: ANOMALOUS



To get more details on channel history, users can go directly to the same view in visualization tools like SmapVitals

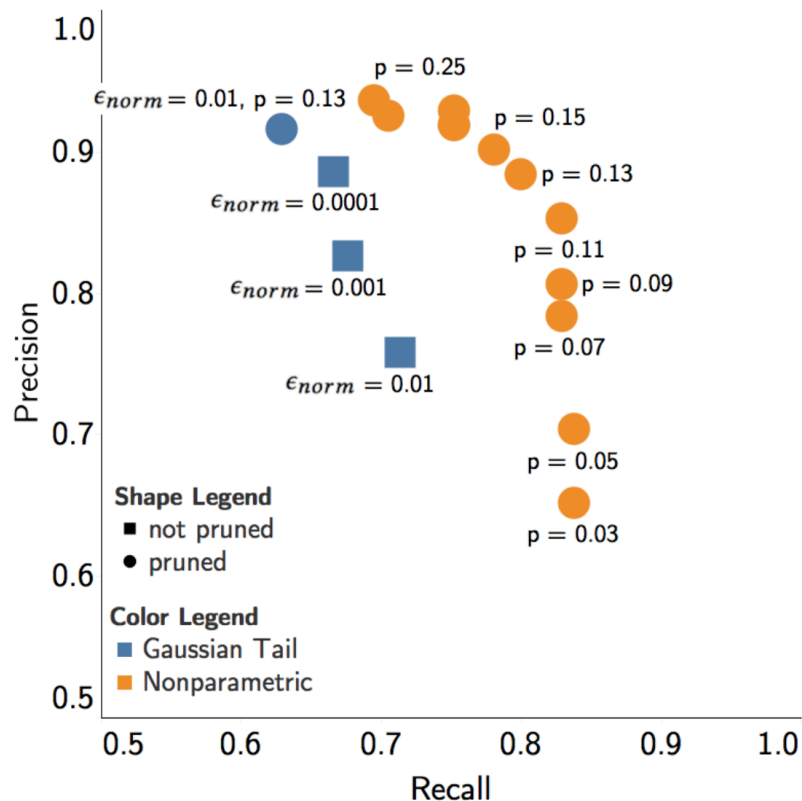
Users can drill down into the raw telemetry for each channel (blue) and compare to the model predictions (orange)

Users can click to tag anomalies as true or false positives, which are used by the system to refine results

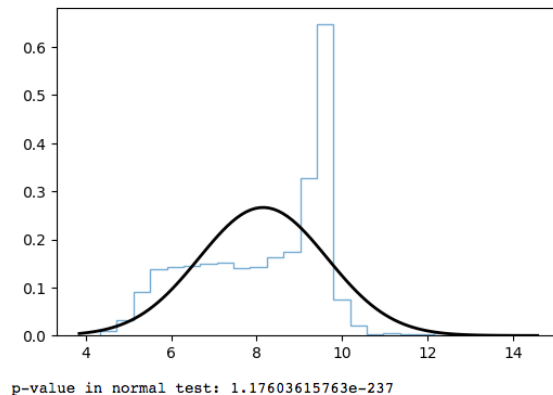
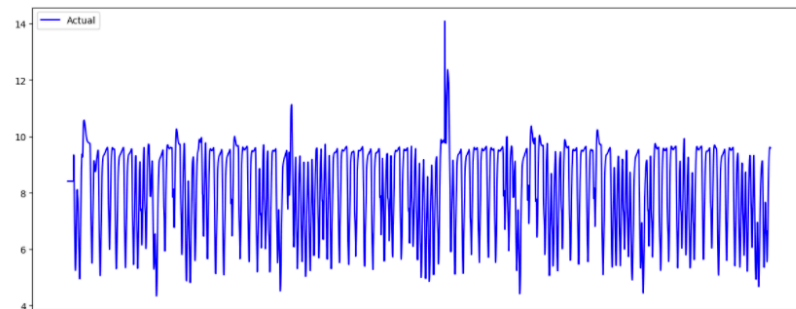
True = green  
False = gray  
Unlabeled = red

Where prediction errors are large, anomalies are flagged

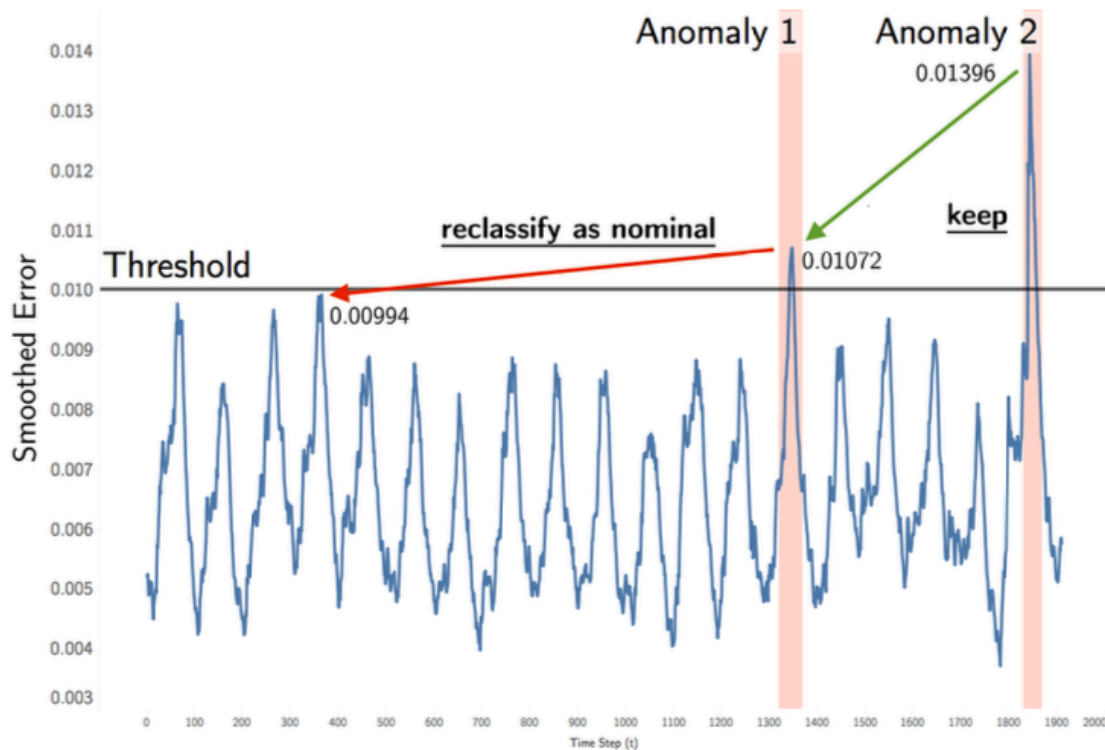
# Results



## Not Gaussian



# Pruning



$$\mathbf{e}_{max} = [0.01396, 0.01072, 0.00994]$$

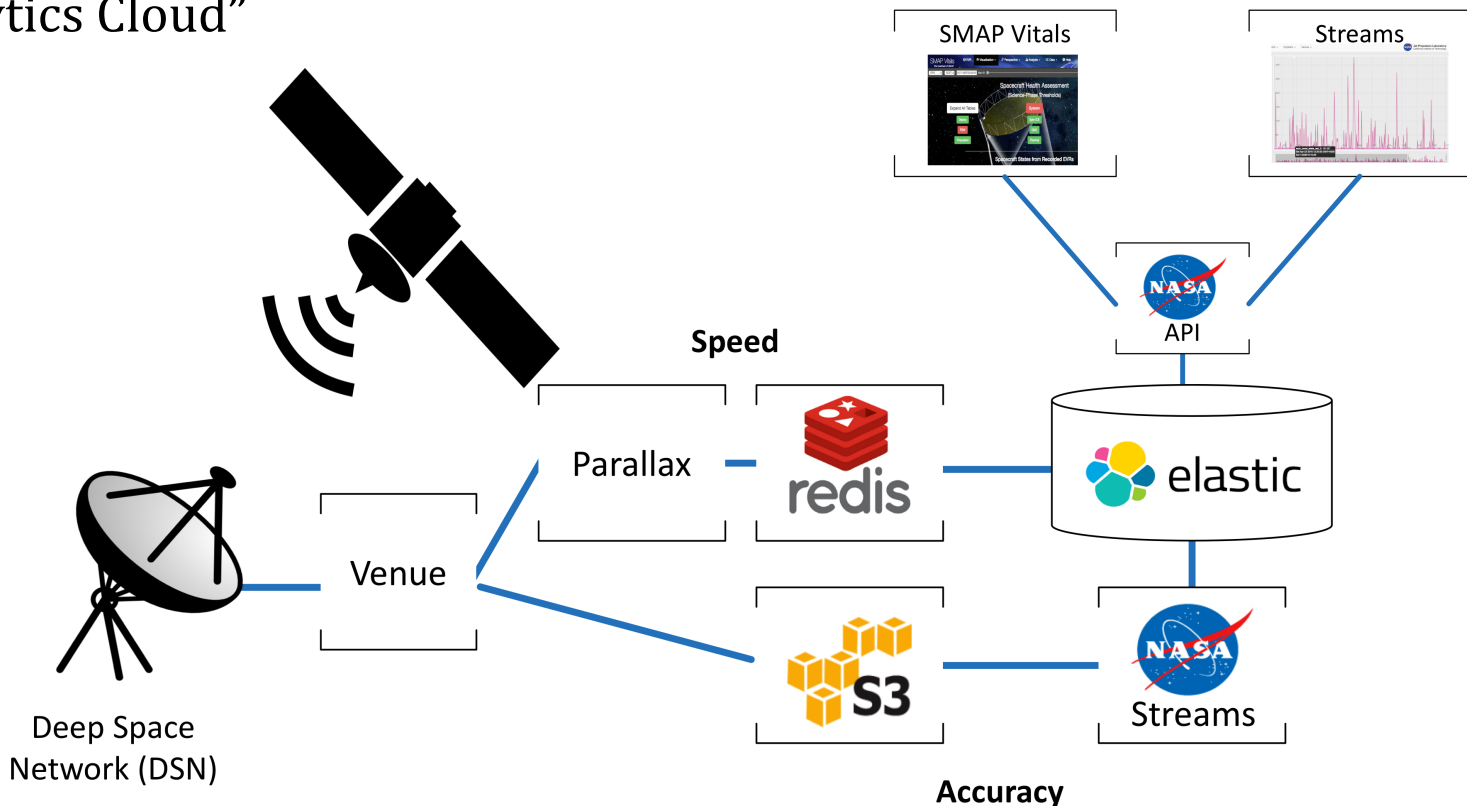
$$p = 0.1$$

Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.



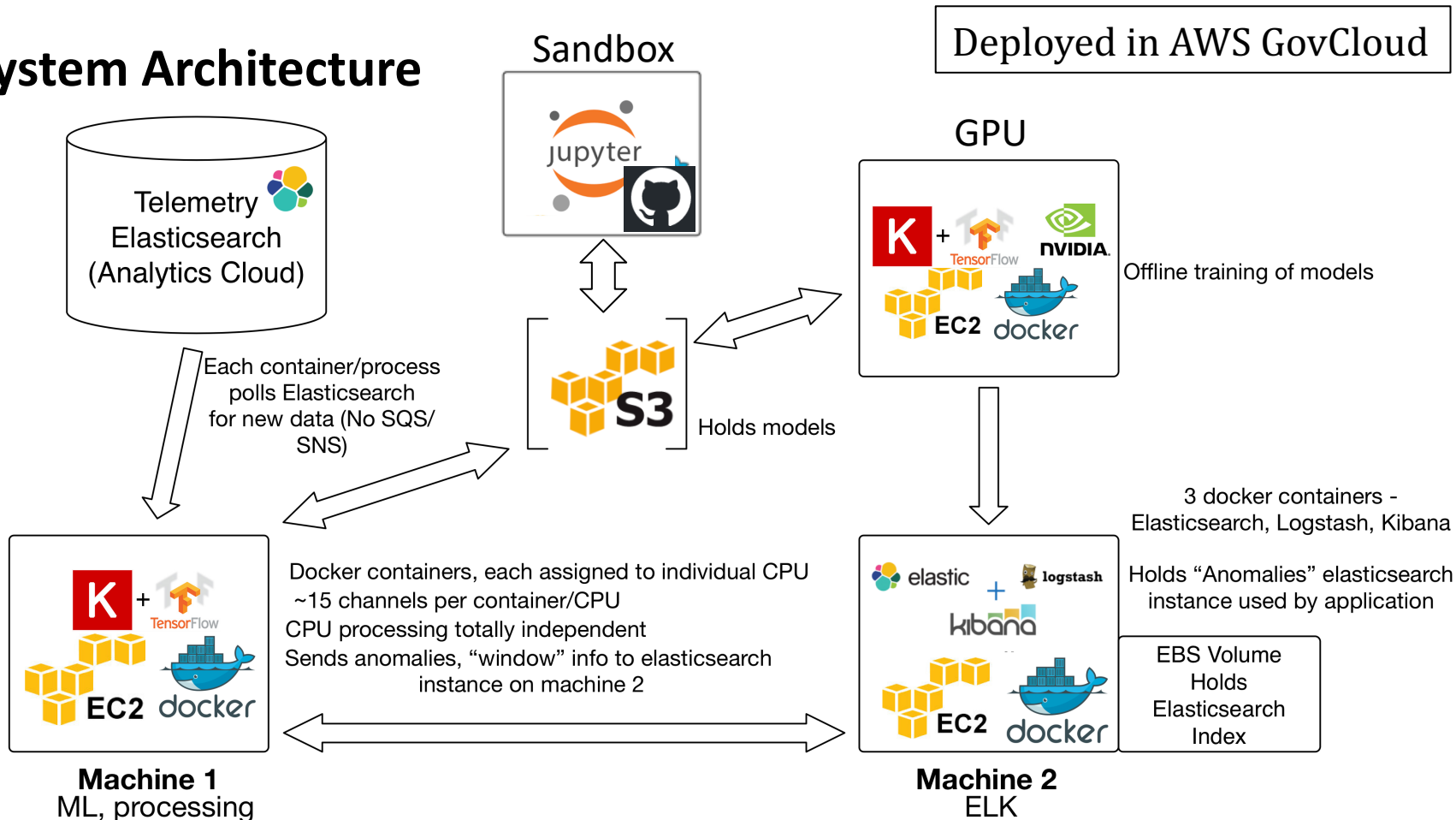
# Foundation

## The “Analytics Cloud”



Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not constitute or imply its endorsement by the United States Government or the Jet Propulsion Laboratory, California Institute of Technology.

# System Architecture



# Soil Moisture Active Passive (SMAP)

- Routine operations
- Major radar failure
- ~4,000 telemetry channels
  - Power, CPU, RAM, Thermal, Radiation, counters
  - 14 command modules
  - 4B values
- Challenges
  - Semi-supervised
  - Complexity, diversity
  - Scale

